

Statistical Diagnostics for Nonlinear Models with Right-Censored Data

Zhong Cheng, Yu Feng

School of Science, Nanjing University of Science and Technology, Nanjing Jiangsu
Email: 396840578@qq.com

Received: Nov. 20th, 2018; accepted: Dec. 6th, 2018; published: Dec. 13th, 2018

Abstract

This paper considers how to solve statistical diagnosis problem of the nonlinear models with right-censored data. First, we use the method of maximum likelihood estimates to reach the parameters. Based on the idea of case-deletion models, we obtain the formula of the parameters and propose some diagnostic statistics to determine outliers or influential points. Finally, we use numerical simulation analysis to verify the feasibility of the theory.

Keywords

Nonlinear Model, Right-Censored Data, Statistical Diagnostics, Case-Deletion Model, Generalized Cook Distance

非线性模型在右删失数据下的统计诊断

程忠, 冯予

南京理工大学, 理学院, 江苏 南京
Email: 396840578@qq.com

收稿日期: 2018年11月20日; 录用日期: 2018年12月6日; 发布日期: 2018年12月13日

摘要

本文研究了非线性模型在带有右删失数据下如何统计诊断的问题。首先用极大似然估计的方法求出了参数的估计问题。再对模型运用数据删除思想进行考量, 求出了删除前后参数的公式, 给出了判定影响点和异常点的一些统计量。最后, 通过数值模拟分析, 证明了该方法的可行性。

关键词

非线性模型, 右删失数据, 统计诊断, 数据删除模型, 广义Cook距离

Copyright © 2018 by authors and Hans Publishers Inc.

This work is licensed under the Creative Commons Attribution International License (CC BY).

<http://creativecommons.org/licenses/by/4.0/>



Open Access

1. 引言

统计诊断[1]是数据分析的重要组成部分, 其主要任务就是通过诊断统计量检测已知观测数据在既定模型拟合时的合理性, 数据删除法[2]是其最基础的方法, 即考虑删除前后统计量的变化。非线性模型[3]可以结合不同的数据类型进行研究, 包括纵向数据、删失数据[4]等, 最近一段时机在纵向数据下运用比较广, 但关于删失数据则很少。右删失数据也是生存分析[5]尤其是在寿命观测中很常见的一种数据。本文考虑了带右删失数据下非线性模型的统计诊断问题[6], 有一定的理论和实践价值。

2. 右删失数据下非线性模型的参数估计

2.1. 右删失数据下的非线性模型

给定观测数据 $(x_i, y_i), i=1, 2, \dots, n$, 非线性回归模型表示为 $y_i = f(x_i, \beta) + \varepsilon_i, i=1, \dots, n$, 其中 $f(x_i, \beta)$ 为参数 β 的非线性函数, $\beta = (\beta_1, \dots, \beta_p)^T$ 为未知的回归系数向量, ε_i 为随机误差。通常假定序列 $\varepsilon_1, \dots, \varepsilon_n$ 独立正态同分布, 其均值为 0, 方差为 σ^2 。该模型的矩阵形式可表示为 $Y = f(X, \beta) + \varepsilon = f(\beta) + \varepsilon$, 其中 $Y = (y_1, \dots, y_n)^T$, $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)^T$, $f(X, \beta) = f(\beta) = (f(x_1, \beta), \dots, f(x_n, \beta))^T$ 。

因为 y_i 右删失, 不失一般性, 我们考虑前 r 个为由于试验的终止而未寿终数据, 即 y_1, \dots, y_r 为右删失数据。由于有删失数据存在, 因此需要求出新的似然函数来进行参数估计。

2.2. 右删失数据下的似然函数

设 X_1, X_2, \dots, X_n 是来自分布 G 的随机变量, 并且独立同分布, 其概率密度函数为 $g(x, \theta)$, θ 为模型参数。同时假设右删失时间 $M_i (i=1, 2, \dots, n)$, 分布为 F 。假设 X_i 和 M_i 相互独立, 记 $Y_i = \min(X_i, M_i)$, $\delta_i = I(M_i \leq X_i)$, 实际观察样本为 (Y_i, δ_i) , 则右删失数据下的似然函数[7]为

$$L((X_1, \delta_1), \dots, (X_n, \delta_n)) = \prod_{i=1}^n [g(X_i, \theta)]^{\delta_i} [1 - G_\theta(X_i)]^{1 - \delta_i} \quad (1)$$

2.3. 模型的极大似然估计

令 $\phi(y) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{y^2}{2}\right)$, 它是标准正态分布的概率密度函数。记 $\beta(y) = \int_y^{+\infty} \phi(t) dt$, $t_i = \frac{y_i - f(x_i, \beta)}{\sigma}$, $i=1, \dots, n$, $S(y) = \frac{\phi(y)}{\Phi(y)}$ 。则根据(1)式, 模型的联合似然函数可表示为 $L = \frac{1}{\sigma^{n-r}} \prod_{i=1}^r \Phi(t_i) \prod_{i=r+1}^n \phi(t_i)$, 取对数, 得对数似然函数为

$$l(\beta) = \sum_{i=1}^r \log \Phi(t_i) - \frac{n-r}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=r+1}^n (y_i - f(x_i, \beta))^2 \quad (2)$$

$$\begin{aligned}
\frac{\partial^2 l(\beta)}{\partial \beta \partial \beta^T} &= \frac{1}{\sigma} \sum_{i=1}^r \frac{\partial S(t_i)}{\partial \beta} \frac{\partial f(x_i, \beta)}{\partial \beta} + \frac{1}{\sigma} \sum_{i=1}^r S(t_i) \frac{\partial f(x_i, \beta)}{\partial \beta \partial \beta^T} \\
&\quad - \frac{1}{\sigma^2} \sum_{i=r+1}^n \left(\frac{\partial f(x_i, \beta)}{\partial \beta} \right)^2 + \frac{1}{\sigma^2} \sum_{i=r+1}^n (y_i - f(x_i, \beta)) \frac{\partial f(x_i, \beta)}{\partial \beta \partial \beta^T} \\
&= -\frac{1}{\sigma^2} \sum_{i=1}^r \left(\frac{\partial f(x_i, \beta)}{\partial \beta} \right)^2 \frac{\phi'(t_i) \Phi(t_i) + \phi^2(t_i)}{\Phi^2(t_i)} \\
&\quad - \frac{1}{\sigma^2} \sum_{i=r+1}^n \left(\frac{\partial f(x_i, \beta)}{\partial \beta} \right)^2 + \frac{1}{\sigma^2} \left\{ \sum_{i=1}^n \frac{\partial^2 f(x_i, \beta)}{\partial \beta \partial \beta^T} (z_i - f(x_i, \beta)) \right\} \\
&= -\frac{1}{\sigma^2} D^T(\beta) \Omega^{-1} D(\beta) + \frac{1}{\sigma^2} [e^T][W]
\end{aligned}$$

即得(5)式。

由于 $E[-\dot{l}(\beta)] = I(\beta)$, 因此可得(6), 证毕。

将以上结果代入高斯-牛顿迭代公式 $\beta^{i+1} = \beta^i + [-\dot{l}(\beta^i)]^{-1} \dot{l}(\beta^i)$ 可得

$$\beta^{i+1} = \beta^i + [D^T(\beta) \Omega^{-1} D(\beta)]^{-1} D^T(\beta) e(\beta)$$

3. 模型诊断

3.1. 数据删除模型

数据删除模型是最基本的统计诊断模型, 比较第 i 个点 (x_i, y_i) 删除前后估计量或统计量之间差异, 其中 i 的取值范围为 $(r+1 \leq i \leq n)$ 。其中, 本文只考虑删除完整的数据部分, 对于由于右删失得到的数据点不研究异常或强影响点问题[8]。

于是, 右删失下非线性模型的数据删除模型可表示为 $y_j = f(x_j, \beta) + \varepsilon_j, j \neq i, r+1 \leq i \leq n$, 其中 y_1, \dots, y_r 为右删失数据, y_{r+1}, \dots, y_n 为完整值。其相应的极大似然估计记为 $\hat{\beta}(i)$ 。

通常可由 $\hat{\beta}(i)$ 的一阶近似公式来比较数据删除前 $\hat{\beta}$ 和删除后 $\hat{\beta}(i)$ 的差异, 即

$$\hat{\beta}^i(i) = \hat{\beta} + [I_{(i)}(\hat{\beta})]^{-1} \dot{l}_{(i)}(\hat{\beta}) \quad (7)$$

其中 $I_{(i)}(\beta)$ 和 $\Omega_{(i)}^{-1}$ 是 $I(\beta)$ 和 Ω^{-1} 删除第 i 点以后得到的矩阵, 而 $\dot{l}_{(i)}(\beta)$ 是 $\dot{l}(\beta)$ 删除第 i 点以后的向量。由引理 1 可知:

$$\dot{l}_{(i)} = \frac{1}{\sigma^2} D_{(i)}^T(\beta) e_{(i)}(\beta) \quad (8)$$

$$I_{(i)}(\beta) = \frac{1}{\sigma^2} D_{(i)}^T(\beta) \Omega_{(i)}^{-1} D_{(i)}(\beta) \quad (9)$$

其中 $D_{(i)}(\beta)$ 为 $D(\beta)$ 删除第 i 行以后的 $(n-1) \times p$ 矩阵, $e_{(i)}(\beta)$ 为 $e(\beta) = Z - f(\beta)$ 删除第 i 点以后的 $(n-1)$ 维向量。

定理 对于右删失下非线性模型, $\hat{\beta}(i)$ 的一阶近似可表示为

$$\hat{\beta}^i(i) = \hat{\beta} - \left\{ \frac{(D^T \Omega^{-1} D)^{-1} d_i e_i}{1 - h_{ii}} \right\}_{\hat{\beta}} \quad (10)$$

其中, d_i 表示 D 的第 i 个行组成的 p 维向量, h_{ii} 为矩阵 $H = D(D^T \Omega^{-1} D)^{-1} D^T$ 的第 i 个对角元。

证明: 将(8)、(9)式代入(7)式可得

$$\hat{\beta}^l(i) = \hat{\beta} + \left\{ \left[D_{(i)}^T(\beta) \Omega_{(i)}^{-1} D_{(i)}(\beta) \right]^{-1} D_{(i)}^T(\beta) e_{(i)}(\beta) \right\}_{\hat{\beta}}$$

其中 $D^T \Omega^{-1} D$ 可进行以下分解:

$$D^T \Omega^{-1} D = \sum_{k=1}^n v_k d_k d_k^T = \sum_{k \neq i} v_k d_k d_k^T + d_i d_i^T = D_{(i)}^T \Omega_{(i)}^{-1} D_{(i)} + d_i d_i^T$$

因此有 $D_{(i)}^T \Omega_{(i)}^{-1} D_{(i)} = D^T \Omega^{-1} D - d_i d_i^T$, 同理可得 $D^T e = \sum_{k \neq i} d_k e_k + d_i e_i = D_{(i)}^T e_{(i)} + d_i e_i$ 。

这些公式代入以上 $\hat{\beta}^l(i)$ 的表达式可得

$$\hat{\beta}^l(i) = \hat{\beta} + \left\{ \left[D^T \Omega^{-1} D - d_i d_i^T \right]^{-1} \left[D^T e - d_i e_i \right] \right\}_{\hat{\beta}}$$

应用矩阵和式求逆公式 $(A + MN)^{-1} = A^{-1} - A^{-1} M (I + N A^{-1} M)^{-1} N A^{-1}$,

取 $A = D^T \Omega^{-1} D, M = -d_i, N = d_i^T$, 所以

$$\begin{aligned} & \left[D^T \Omega^{-1} D - d_i d_i^T \right]^{-1} \\ &= \left(D^T \Omega^{-1} D \right)^{-1} - \left(D^T \Omega^{-1} D \right)^{-1} (-d_i) \left[I + d_i^T \left(D^T \Omega^{-1} D \right)^{-1} (-d_i) \right]^{-1} d_i^T \left(D^T \Omega^{-1} D \right)^{-1} \\ &= \left(D^T \Omega^{-1} D \right)^{-1} + \frac{\left(D^T \Omega^{-1} D \right)^{-1} d_i d_i^T \left(D^T \Omega^{-1} D \right)^{-1}}{1 - h_{ii}} \end{aligned}$$

由于 $i(\hat{\beta}) = \frac{1}{\sigma^2} D^T(\hat{\beta}) e(\hat{\beta}) = 0$ 。

因此, 当 $\beta = \hat{\beta}$ 时,

$$\begin{aligned} \left[D^T \Omega^{-1} D - d_i d_i^T \right]^{-1} \left[D^T e - d_i e_i \right] &= - \left(D^T \Omega^{-1} D \right)^{-1} d_i e_i - \frac{\left(D^T \Omega^{-1} D \right)^{-1} d_i d_i^T \left(D^T \Omega^{-1} D \right)^{-1} d_i e_i}{1 - h_{ii}} \\ &= - \left(D^T \Omega^{-1} D \right)^{-1} d_i e_i \left[1 + \frac{h_{ii}}{1 - h_{ii}} \right] = - \frac{\left(D^T \Omega^{-1} D \right)^{-1} d_i e_i}{1 - h_{ii}} \end{aligned}$$

所以,

$$\hat{\beta}^l(i) = \hat{\beta} - \left\{ \frac{\left(D^T \Omega^{-1} D \right)^{-1} d_i e_i}{1 - h_{ii}} \right\}_{\hat{\beta}}$$

证明完成。

以上定理理解出来了数据删除前和删除后估计量之间的数学公式, 由此公式出发还能得到许多其他的估计量。这个公式表明如果 $\hat{\beta}^l(i)$ 与 $\hat{\beta}$ 之间应相差很大, 则说明 i 可能为异常点或强影响点。

3.2. 统计诊断量

3.2.1. 广义 Cook 距离

根据公式(10), 我们即可得到 $\hat{\beta} - \hat{\beta}^l(i)$ 。在一般的回归模型中, Cook 已经提出了 Cook 距离, 因此这里我们将 Cook 距离推广到带删失数据的非线性模型中。广义 Cook 距离定义为

$$GD_i = \|\hat{\beta} - \hat{\beta}(i)\|_M^2 = \frac{(\hat{\beta} - \hat{\beta}(i))^T M (\hat{\beta} - \hat{\beta}(i))}{c}$$

其中, M 为正定矩阵; $c > 0$ 为尺度因子。此处取为

$$M = D^T \Omega^{-1} D, c = p \tilde{\sigma}^2$$

则可得

$$GD_i = \frac{(\hat{\beta} - \hat{\beta}(i))^T D^T \Omega^{-1} D (\hat{\beta} - \hat{\beta}(i))}{p \tilde{\sigma}^2}$$

若用 $\hat{\beta}'(i)$ 代替 $\hat{\beta}(i)$, 则得到 1 阶近似公式, 将式(10)代入可得广义 Cook 距离:

$$GD_i = \frac{\hat{h}_{ii}}{1 - \hat{h}_{ii}} \frac{r_i^2}{p}, r_i = \frac{\hat{e}_i}{\tilde{\sigma} \sqrt{1 - \hat{h}_{ii}}}$$

3.2.2. 似然距离

似然距离在一般情况下没有似然解, 因此本文采用其一阶近似公式

$$LD_i'(\beta) = (\hat{\beta} - \hat{\beta}(i))^T [I(\hat{\beta})] (\hat{\beta} - \hat{\beta}(i))$$

把 $\hat{\beta}'(i)$ 换成 $\hat{\beta}(i)$, 即可得似然距离的近似解

$$LD_i'(\beta) = (\hat{\beta} - \hat{\beta}'(i))^T [I(\hat{\beta})] (\hat{\beta} - \hat{\beta}'(i)) = pGD_i'$$

因此能得知, 广义 Cook 距离与似然距离只多出了 p 。

4. 实例分析

下面用一个有关渔场捕鱼的数据来进行数值模拟分析, 验证非线性模型在带有带右删失数据下的统计诊断具有可行性。

鱼卵数量与可捕获的成鱼数量之间的关系, 是经营渔场者十分关心的问题。为研究二者之间的关系, Ricker 和 Smith(1975)给出了在 Skeener 河中红鳟鲑鱼的产卵量和可捕获的成鱼量的测量数据。表 1 [9]列出了这些数据; 表中 x 为鱼卵量, y 为可捕获成鱼的数量。所选用的模型为

$$\log(y_i) = \beta_1 \log(x_i) \exp\{-\beta_2 \log(x_i)\} + \varepsilon_i, \varepsilon_i \sim N(0, \sigma^2), i = 1, 2, \dots, 28$$

考虑到 y 的前 3 个数据删失, 假设 $y_i > 500 (1 \leq i \leq 3)$ 。数据如表 1 所示。

Table 1. Data of red trout salmon

表 1. 红鳟鲑鱼数据

i	YEAR	X	Y
1	1940	963	500
2	1941	572	500
3	1942	305	500
4	1943	272	438
5	1944	824	3071
6	1945	940	957
7	1946	486	934

Continued

8	1947	307	971
9	1948	1066	2257
10	1949	480	1451
11	1950	393	686
12	1951	176	127
13	1952	237	700
14	1953	700	1381
15	1954	511	1393
16	1955	87	363
17	1956	370	668
18	1957	448	2067
19	1958	819	644
20	1959	799	1747
21	1960	273	744
22	1961	936	1087
23	1962	558	1335
24	1963	597	1981
25	1964	848	627
26	1965	619	1099
27	1966	397	1532
28	1967	616	2086

4.1. 参数估计

根据以上数据表, 我们能解出参数 β 的估计值:

$$\hat{\beta} = (1.5864, 0.0576)^T$$

4.2. 影响分析

估计完参数之后, 由数据删除模型计算得到广义 Cook 距离。图 1 是广义 Cook 距离的散点图。

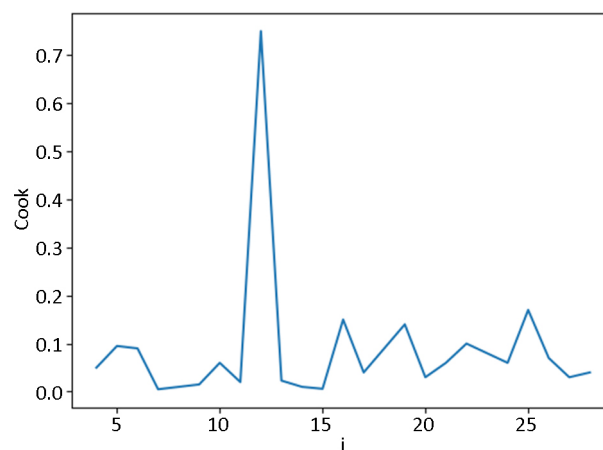


Figure 1. Scatter plot of generalized Cook distance
图 1. 广义 Cook 距离的散点图

从图 1 可以看出: 第 12 号点异于其他点, 而且不涉及删失部分, 因此第 12 号点为异常点。

5. 结束语

本文介绍了非线性模型在带有右删失数据下如何进行统计诊断的问题, 给出了参数估计的方法以及如何通过统计诊断量判断强影响点或异常点, 都是很经典很实际的办法; 最后通过数值模拟分析, 验证了诊断方法可行性。考虑到经典方法的局限性, 本模型还可以用经验似然或其他方法进行更深一步的研究, 这也是本文作者努力的方向。

基金项目

国家自然科学基金资助项目(11271189)。

参考文献

- [1] 韦博成, 鲁国斌, 等. 统计诊断引论[M]. 南京: 东南大学出版社, 1991.
- [2] 翟爽. 基于数据删除的广义线性模型诊断方法[D]: [硕士学位论文]. 哈尔滨: 东北林业大学理学院, 2012.
- [3] 胡宏昌, 崔恒建, 秦永松, 等. 近代线性回归分析方法[M]. 北京: 科学出版社, 2013.
- [4] 周勇. 广义估计方程估计方法[M]. 北京: 科学出版社, 2013.
- [5] 陈家鼎. 生存分析与可靠性[M]. 北京: 北京大学出版社, 2005.
- [6] 王思洋, 胡涛, 崔恒建. 删失数据非线性回归模型的广义 M 估计[J]. 北京师范大学学报(自然科学版), 2014, 50(1): 1-6.
- [7] 朱成莲. 带右删失数据的非线性模型的参数估计[J]. 统计与决策, 2009(14): 155-156.
- [8] 季文奇, 冯予. 右删失数据下广义线性模型的统计诊断[J]. 重庆理工大学学报(自然科学), 2017, 31(8): 174-181.
- [9] 韦博成, 林金官, 解锋昌. 统计诊断[M]. 北京: 高等教育出版社, 2009.

知网检索的两种方式:

1. 打开知网页面 <http://kns.cnki.net/kns/brief/result.aspx?dbPrefix=WWJD>
下拉列表框选择: [ISSN], 输入期刊 ISSN: 2325-2251, 即可查询
2. 打开知网首页 <http://cnki.net/>
左侧“国际文献总库”进入, 输入文章标题, 即可查询

投稿请点击: <http://www.hanspub.org/Submission.aspx>

期刊邮箱: sa@hanspub.org